

### 5.2.6. Executie neordonata si redenumirea registrelor

Numeroase UCP-uri moderne sunt atât în banda de asamblare cât și superscalare, așa cum se arată în Fig. 5.27. Ceea ce înseamnă în general, că există o unitate de citire care citește instrucțiuni din memorie înainte ca acestea să fie necesare pentru a fi trimise unității de decodificare. Unitate de decodificare trimite instrucțiunile decodificate unităților functionale adecvate pentru execuție. În anumite cazuri, instrucțiuni individuale se pot sparge în micro-operatii (micro-ops) înainte de a le trimite unităților functionale depinzând de ce sunt capabile să facă unitățile functionale.

Evident, proiectarea masinilor este mai simplă dacă toate instrucțiunile sunt executate în ordinea în care sunt citite (presupunând pentru moment că algoritmul de predicție a ramificațiilor nu greșeste niciodată). Totuși, execuția în ordine nu furnizează întotdeauna performanță optimă din cauza dependenței între instrucțiuni. Dacă o instrucțiune necesită o valoare calculată de către instrucțiunea precedentă, cea de-a doua instrucțiune nu își poate începe execuția până ce prima nu a produs valoarea necesară. În această situație (o dependență RAW) a doua instrucțiune trebuie să aștepte. Există de asemenea și alte tipuri de dependență, așa cum vom vedea în curând.

În încercarea de a aborda aceste probleme și de a realiza o performanță superioară, anumite UCP-uri permit sărirea unor instrucțiuni pentru a ajunge la instrucțiuni următoare care nu sunt dependente. Nu mai trebuie să spunem că algoritmul intern de planificare a instrucțiunilor utilizat trebuie să furnizeze același efect ca și când programul ar fi executat în ordinea scrisă. Vom demonstra acum cum funcționează reordonarea instrucțiunilor utilizând un exemplu detaliat.

Pentru a ilustra natura problemei vom începe cu o mașină care întotdeauna furnizează instrucțiuni în ordinea programului și de asemenea necesită terminarea execuției în ordinea programului. Semnificația ideii va deveni clară mai târziu.

Mașina din exemplul nostru are opt registre vizibile programatorului, de la R0 la R7. Toate instrucțiunile aritmetice utilizează trei registre: două pentru operanți și unul pentru rezultat. Vom presupune că dacă o instrucțiune este decodificată în ciclul  $n$ , execuția începe în ciclul  $n+1$ . Pentru o instrucțiune simplă, cum ar fi o adunare sau o înmulțire, scrierea în registrul destinație se produce la sfârșitul ciclului  $n+2$ . Pentru o instrucțiune mai complicată, cum ar fi o înmulțire, scrierea se produce la sfârșitul ciclului  $n+3$ . Pentru a apropia exemplul de realitate, vom permite unității de decodificare să lanseze în execuție până la două instrucțiuni pe ciclu de ceas. UCP-urile superscalare comerciale adesea pot lansa în execuție patru sau chiar șase instrucțiuni pe ciclu de ceas.

Secvența de execuție din exemplul nostru este arată în Fig. 5.31. Aici prima coloană furnizează numărul de cicluri și a doua coloană furnizează numărul instrucțiunii. A treia coloană listează instrucțiunea decodificată, a patra spune care instrucțiune este lansată în execuție (cu un maxim de două pe ciclu de ceas). A cincea spune care instrucțiune a fost retrasă (terminată). Să ne amintim că în acest exemplu cerem atât lansarea în execuție în ordine cât și terminarea în ordine, astfel ca instrucțiunea  $k+1$  nu poate fi lansată în execuție până când instrucțiunea  $k$  nu a fost lansată în execuție iar instrucțiunea  $k+1$  nu poate fi retrasă (înțelegând efectuarea scrierii în registrul destinație) până ce instrucțiunea  $k$  nu a fost retrasă. Celelalte 16 coloane sunt discutate mai jos.

C	#	Decodificat	Iss	Ret	Registre care sunt citite								Registre care sunt scrise							
					0	1	2	3	4	5	6	7	0	1	2	3	4	5	6	7
1	1	R3=R0*R1	1		1	1										1				
	2	R4=R0+R2	2		2	1	1									1	1	1		
2	3	R5=R0+R1	3		3	2	1									1	1	1		
	4	R6=R1+R4	-		3	2	1									1	1	1		
3					3	2	1									1	1	1		
4				1	2	1	1										1	1		
				2	1	1												1		
				3																
5			4			1			1										1	
	5	R7=R1*R2	5			2	1		1										1	1
6	6	R1=R0-R2				2	1		1										1	1
7				4		1	1													1
8				5																
9			6		1		1							1						
	7	R3=R3*R1	7		1	1	1	1						1		1				
10					1	1	1	1						1		1				
11				6		1		1												
12				7																
13		R1=R4+R4	8						2					1						
14									2					1						
15				8																

C = Ciclu; Iss = Lansata; Ret = Retrasa

Figura 5.31. Functionarea unei UCP superscalar cu lansare în executie în ordine si terminare în ordine.

Dupa decodificarea unei instructiuni unitatea de decodificare trebuie sa decida daca poate fi lansata în executie imediat sau nu. Pentru a lua aceasta decizie, unitatea de decodificare trebuie sa cunoasca starea tuturor registrelor. Daca, de exemplu, instructiunea curenta necesita un registru a carui valoare nu a fost încă calculata, instructiunea curenta nu poate fi lansata în executie si UCP trebuie sa se blocheze.

Vom pastra informatii despre utilizarea registrelor folosind un dispozitiv numit tabela de scor (scoreboard), care a fost prima data prezent în CDC 6600. Tabela de scor are un mic numarator pentru fiecare registru indicând de câte ori este utilizat acel registru ca sursa de catre instructiunile aflate în executie. Daca maximum, sa zicem, 15 instructiuni se pot executa simultan atunci este necesar un contor de 4 biti. Când este lansata în executie o instructiune, intrările tabelii de scor pentru registrele sale cu operanzi sunt incrementate. Când o instructiune este retrasa, intrările sunt decrementate.

Tabela de scor are de asemenea un numarator pentru fiecare registru, pentru a pastra informatii despre registrele utilizate ca destinatii. Deoarece se permite o singura scriere la un moment dat, aceste numaratoare pot fi de 1 bit. În Fig. 5.31 cele 16 coloane din dreapta arata tabela de scor.

În masinile reale, tabela de scor pastreaza si informatii despre utilizarea unitatilor functionale pentru a împiedica lansarea în executie unei instructiuni pentru care nu este disponibila nici o unitate functionala. Pentru simplitate, vom presupune ca întotdeauna exista o unitate functionala adecvata disponibila, asa ca nu vom arata unitatile functionale pe tabeta de scor.

Prima linie din Fig. 5.31 arata I1 (instructiunea 1), care înmulteste R0 cu R1 si pune rezultatul în R3. Deoarece nici unul din aceste registre nu este deocamdata folosit, instructiunea este lansata în executie si tabela de scor actualizata pentru a reflecta faptul ca R0 si R1 sunt citite si R3 este înscris. Nici una din instructiunile urmatoare nu poate scrie în oricare din aceste registre sau nu poate citi R3 pâna când I1 nu a fost retrasa. Deoarece aceasta instructiune este o înmultire, ea se va încheia la sfârșitul ciclului 4. Valorile tabelii de scor aratate pe fiecare linie reflecta starea acestora dupa ce instructiunea de pe acea linie a fost lansata în executie. Intrările vide (nemarcate) sunt zerouri.

Deoarece exemplul nostru este o masina superscalara care poate lansa în executie doua instructiuni pe ciclu, o a doua instructiune (I2) este lansata în executie în timpul ciclului 1. Aceasta aduna R0 si R2, memorând rezultatul în R4. Pentru a vedea daca aceasta instructiune poate fi lansata în executie se aplica urmatoarele reguli:

1. Daca orice operand este în curs de înscriere, nu se va lansa în executie (dependenta RAW).
2. Daca registrul rezultat este în curs de citire, nu se va lansa în executie (dependenta WAR).
3. Daca registrul rezultat este în curs de scriere, nu se va lansa în executie (dependenla WAW).

Am discutat deja dependentele RAW, care au loc când o instructiune trebuie sa utilizeze ca sursa un rezultat pe care o instructiune precedenta încă nu l-a produs. Celelalte doua dependente sunt mai putin serioase.

Ele reprezintă de fapt conflicte de resurse. Într-o dependență WAR (Write After Read, rom: scriere după citire), o instrucțiune încearcă să supraînregistre un registru pe care o instrucțiune precedentă poate nu a terminat încă să-l citească. O dependență WAW (Write After Write, TOJ scriere după scriere) este similară. Acestea pot fi adesea evitate prin plasarea rezultatelor celei de-a doua instrucțiuni undeva în alta parte (eventual temporar). Dacă nici una din cele trei dependente de mai sus nu există și dacă unitatea funcțională necesară este disponibilă, instrucțiunea este lansată în execuție. În acest caz, I2 utilizează un registru (R0) care este în curs de citire de către o instrucțiune în desfășurare, dar această suprapunere este permisă deci I2 este lansată în execuție. Asemănător, I3 este furnizată în timpul ciclului 2. Trecem acum la I4, care necesită utilizarea lui R4. Din nefericire, vedem din linia 3 că R4 este în curs de scriere. Aici avem o dependență RAW, așa că unitatea de decodificare se blochează până când R4 devine disponibil. Cât timp este blocată, aceasta nu va prelua instrucțiuni de la unitatea de citire a instrucțiunilor. Când tampoanele interne ale unității de citire se încarcă integral, se va opri citirea instrucțiunii în avans.

Merita să observăm că următoarea instrucțiune în ordinea programului, I5 nu are conflicte cu nici una din instrucțiunile în desfășurare. Ar fi putut fi decodificată și lansată în execuție, dar nu este căci acest proiect necesită lansarea instrucțiunilor în execuție în ordine.

Să urmărim acum ce se întâmplă în timpul ciclului 3. I2, fiind o adunare (două cicluri), se încheie la sfârșitul ciclului 3. Din nefericire, ea nu poate fi retrasă (eliberând astfel R4 pentru I4). De ce nu? Motivul este că acest proiect necesită, de asemenea, retragerea în ordine. De ce? Ce prejudiciu posibil ar putea proveni din memorarea în R4 acum și marcarea acestuia ca disponibil?

Răspunsul este subtil, dar important. Să presupunem că instrucțiunile se pot termina în ne-ordine.

Atunci, dacă apare o întrerupere, ar fi foarte dificil să se salveze starea mașinii, astfel încât să se restabilească ulterior. În particular nu ar fi posibil să se spună că toate instrucțiunile până la o anumită adresă au fost executate și toate instrucțiunile de după aceasta nu. Aceasta este numită întrerupere precisă (precise interrupt) și este o caracteristică dezirabilă într-o UCP (Moudgill și Vassiliadis, 1996). Retragera în ne-ordine face întreruperile imprecise, motiv pentru care unele mașini necesită terminarea în ordine a instrucțiunilor.

Revenind la exemplul nostru la sfârșitul ciclului 4 toate cele trei instrucțiuni în desfășurare pot fi retrase, astfel că în ciclul 5 I4 poate fi în sfârșit lansată în execuție, alături de nou decodificată I5. Ori de câte ori o instrucțiune este retrasă unitatea de decodificare trebuie să verifice dacă există o instrucțiune blocată care poate fi acum lansată în execuție.

În ciclul 6, I6 se blochează pentru că trebuie să scrie în R1, iar R1 este ocupat. Ea este lansată în sfârșit în ciclul 9. Întreaga secvență de opt instrucțiuni consumă pentru terminare 15 cicluri datorită numeroaselor dependente, chiar dacă hardware-ul este capabil să lanseze în execuție două instrucțiuni în fiecare ciclu. Remarcați totuși, citind de sus în jos coloana Lansată din Fig. 5.31, că toate instrucțiunile au fost lansate în execuție în ordine. În același fel, coloana Retrasă arată că ele au fost retrase de asemenea în ordine.

Să considerăm acum o soluție alternativă: execuția în ne-ordine. În această soluție instrucțiunile pot fi lansate în execuție în alta ordine și pot fi retrase de asemenea în alta ordine. Aceeași secvență de opt instrucțiuni este aratăată în Fig. 5.32, numai că acum sunt permise lansarea în execuție și retragerea în ne-ordine.

Prima diferență se produce în ciclul 3. Chiar dacă I4 s-a blocat, putem decodifica și lansa I5 pentru că nu este în conflict cu nici o instrucțiune în desfășurare. Sarirea peste instrucțiuni generează totuși o nouă problemă. Să presupunem că I5 a utilizat un operand calculat de instrucțiunea peste care s-a sarit, I4. Cu tabela curentă de scor, nu am remarcat acest lucru. Ca o consecință trebuie să extindem tabela de scor pentru a reține informații despre memorările făcute de instrucțiunile peste care s-a sarit. Aceasta se poate face adăugând un al doilea plan de biti, un bit pentru fiecare registru, pentru a reține informații despre memorările făcute de instrucțiunile blocate (aceste contoare nu sunt arătate în figura). Acum regula pentru lansarea instrucțiunilor în execuție trebuie extinsă pentru a evita lansarea oricărei instrucțiuni cu un operand planificat să fie modificat de o instrucțiune care a sosit mai înainte, dar a fost sarită.

Acum să ne întoarcem la I6, I7 și I8 din Fig. 5.32. Aici vedem că I6 calculează o valoare în R1 care este utilizată de I7. Totuși, vedem și că valoarea nu este utilizată din nou pentru că I8 supraînregistre R1.

Nu există nici un motiv real să se utilizeze R1 ca loc de păstrare a rezultatului lui I6. Mai rău, R1 este o alegere teribilă de registru intermediar, deși perfect rezonabilă pentru un compilator sau programator obișnuit cu ideea execuției secvențiale fără suprapunere de instrucțiuni.

În Fig. 5.32 introducem o nouă tehnică pentru rezolvarea acestei probleme: redenumirea registrelor (register renaming). Unitatea inteligentă de decodificare schimbă utilizarea lui R1 în I6 (ciclul 3) și I7 (ciclul 4) într-un registru secret, S1, invizibil programatorului. Acum I6 poate fi lansată concurent cu I5. UCP-urile moderne au adesea zeci de registre secrete folosite pentru redenumirea registrelor. Aceasta tehnica poate elimina adesea dependențele WAR și WAW.

C	#	Decodificat	Iss	Ret	Registre care sunt citite								Registre care sunt scrise							
					0	1	2	3	4	5	6	7	0	1	2	3	4	5	6	7
1	1	R3=R0*R1	1		1	1										1				
	2	R4=R0+R2	2		2	1	1									1	1			
2	3	R5=R0+R1	3		3	2	1									1	1	1		
	4	R6=R1+R4	-		3	2	1									1	1	1		
3	5	R7=R1*R2	5		3	3	2									1	1	1		1
	6	S1=R0-R2	6	2	4	3	3									1	1	1		1
					3	3	2									1		1		1
4	7	R3=R3*R1 S2=R4+R4	4 - 8		3	4	2		1							1		1	1	1
					3	4	2		1							1		1	1	1
					3	4	2		3							1		1	1	1
					2	3	2		3								1	1	1	1
5				1	3	2	2		3									1	1	1
					1	2	2		3										1	1
					6									1						
						2	1		3											
6			7			2	1	1	3					1		1			1	1
						1	1	1	2					1		1				1
					4			1	2											
					5			1	2					1		1				
7					6			1								1				
								1												1
8					4			1								1				
9			7																	

C = Ciclu; Iss = Lansata; Ret = Retrasa

**Figura 5.32. Functionarea unei UCP superscalar cu lansare în ne-ordine si terminare în ne-ordine.**

La I8 utilizam din nou redenumirea registrelor. De aceasta data R1 este redenumit în S2 asa ca adunarea poate fi lansata înainte ca R1 sa fie liber, la sfârșitul ciclului 6. Daca apare necesitatea ca de aceasta data rezultatul sa fie chiar în R1, continutul lui S2 poate fi oricând copiat acolo la timp. Chiar mai bine, toate instructiunile urmatoare care au nevoie de el pot avea sursele lor redenumite la registrul unde este într-adevar memorat. În orice caz, adunare trebuie sa înceapa mai devreme în acest fel.

Pe multe masini reale redenumirea este strâns legata de modul în care sunt organizate registrele. Exista multe registre secrete si o tabela care pune în corespondenta registrele vizibile programatorului în registre secrete. Astfel registrul real utilizat pentru, sa spunem, R0 este localizat urmarind intrare 0 a acestei table de corespondenta. În acest fel, nu exista registrul real R0, ci numai o legatura între numele R0 si unul dintre registrele secrete. Aceasta legatura se schimba frecvent în timpul executiei pentru a evita dependentele.

Observati ca în Fig. 5.32 când se citește de sus în jos coloana a patra, instructiunile nu au fost lansate în ordine. Nici nu au fost retrase în ordine. Concluzia acestui exemplu este simpla: utilizând executia în ne-ordine si redenumirea registrelor suntem capabili sa marim viteza calculului printr-un factor apropiat de doi.

### 5.2.7. Executie speculativa

În sectiunea precedenta am introdus conceptul de reordonare a instructiunilor pentru îmbunatatirea performantei. Desi nu am mentionat explicit punctul central a fost reordonarea instructiunilor în cadrul unui singur bloc. Este acum momentul sa privim acest aspect mai îndeaproape.

Programele pentru calculatoare se pot sparge în blocuri de baza (basic blocks), fiecare bloc de baza constând într-o secventa liniara de cod cu un punct de intrare în capatul de sus si o iesire în capatul de jos. Un bloc de baza nu contine nici un fel de structuri de control (de exemplu instructiuni if sau while) astfel încât translatarea sa în limbaj masina nu contine nici un fel de ramificatii. Blocurile de baza sunt conectate prin instructiuni de control.

Un program în aceasta forma poate fi reprezentat ca un graf orientat, asa cum este aratat în Fig. 5.33. Aici calculam suma cuburilor întregilor pari si impari, pâna la o anumita limita si acumulam rezultatele lor în evensum si respectiv, oddsum. În cadrul fiecarui bloc de baza, tehnicile de reordonare din sectiunile precedente functioneaza bine.

Necazul este ca majoritatea blocurilor de baza sunt scurte si nu exista suficient paralelism în acestea pentru ai exploata efectiv. În consecinta pasul urmator este sa se permita reordonarea pentru a intersecta

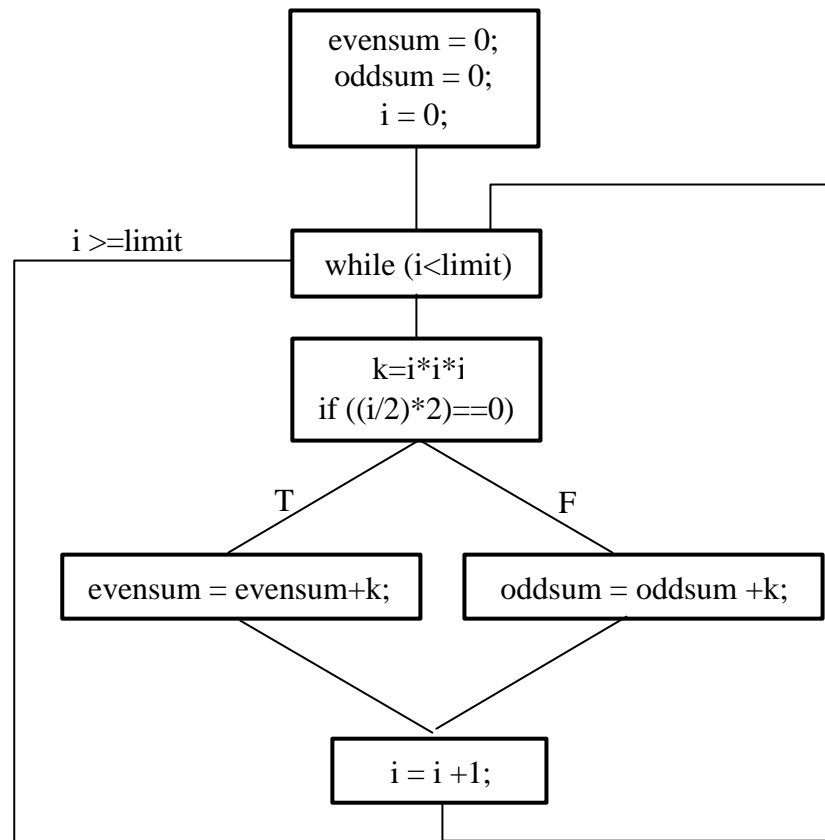
granitele blocurilor de baza în încercarea de a umple toate pozitiile de lansare. Cel mai mare câștig apare când o operație posibil lentă poate fi mutată mai sus în graf pentru a începe mai devreme. Aceasta ar putea să fie instrucțiunea LOAD, o operație în virgula mobilă sau chiar startul unui lanț lung de dependențe. Mutarea de cod în sus peste o ramificație se numește înălțare (hoisting).

```
evensum = 0;
oddsum = 0;
i = 0;
```

```
while (i < limit){
```

```
    k = i * i * i
    if ((i/2)*2==0)
```

```
        evensum = evensum+k
    else
        oddsum = oddsum + k
    i = i + 1
}
```



(a)

(b)

**Fig. 5.33.(a) Un fragment de program. (b) Graful corespunzător blocului**

Să ne imaginăm că în Fig. 5.33 toate variabilele au fost ținute în registre, exceptând `evensum` și `oddsum` (din lipsa de registre). Pare să existe un sens deci să se mute instrucțiunile lor LOAD în vârful buclei, înainte de calcularea lui `k`, pentru a le începe mai devreme, astfel încât valorile să fie disponibile când vor fi necesare. Desigur, numai una din ele va fi necesară la fiecare iterație, așa că celălalt LOAD (va fi irosit, dar, dacă memoria intermediară și memoria principală sunt în banda de asamblare și dacă sunt disponibile poziții de lansare s-ar putea să merite. Executia codului înainte de a ști dacă va fi necesar într-adevăr, se numește executie speculativă (speculative execution). Utilizarea acestei tehnici necesită sprijin din partea compilatorului și a hardware-ului, precum și anumite extensii arhitecturale. În majoritatea cazurilor, reordonarea instrucțiunilor peste granitele blocurilor de bază este mai presus de capacitatea hardware-ului, deci compilatorul trebuie să mute instrucțiunile explicit.

Executia speculativă introduce anumite probleme interesante. Prima ar fi că, este esențial ca una din instrucțiunile speculative să nu existe rezultate irevocabile, pentru că s-ar putea constata mai târziu că n-ar fi trebuit executate. În Fig. 5.33, se poate citi `evensum` și `oddsum` și de asemenea se poate face adunarea imediat ce `k` este disponibil (chiar înainte de instrucțiunea `if`), dar nu este bine să se memoreze rezultatele în memorie. În secvențe de cod mai complicate, o soluție banală de a preveni scrierea registrelor de către codul speculativ, înainte să se știe dacă acest lucru este dorit, este să se redenumască toate registrele destinate utilizate de codul

speculativ. În acest fel, sunt modificate numai registrele temporare, așa ca nu este nici o problema dacă în ultimul moment codul nu mai este necesar. Dacă codul este necesar registrele temporare sunt copiate în adevăratele registre destinație. Așa cum va puteți imagina, menținerea tabelului de marcaj pentru păstrarea tuturor acestor informații nu este simplă, dar se poate face având suficient hardware.

Mai există însă o altă problemă introdusă de codul speculativ, care nu se poate rezolva prin redenumirea registrelor. Ce se întâmplă dacă o instrucțiune executată speculativ cauzează o excepție?

Un exemplu neplăcut, dar nu fatal, este o instrucțiune LOAD care cauzează o rată în memoria intermediară pe o mașină cu o dimensiune mare a liniei de memorie intermediară (să zicem, 256 de octeți) și o memorie cu mult mai lentă decât UCP și memoria intermediară. Dacă un LOAD, care este de fapt necesar, oprește funcționarea mașinii pentru mai multe cicluri până când linia de memorie intermediară este încărcată. Nu este alta soluție din moment ce cuvântul este necesar. Totuși blocarea mașinii pentru citirea unui cuvânt care se dovedește a nu fi necesar este contraproductivă. Prea multe astfel de "optimizări" pot face UCP mai lentă decât dacă nu le-ar avea deloc. (Dacă mașina are memorie virtuală, care este discutată în capitolul 6, un LOAD speculativ poate chiar să genereze o lipsă de pagină, care necesită o operație cu discul pentru a încărca pagina dorită. Evenimentele false lipsă de pagină pot avea un efect teribil asupra performanței, așa ca este important să fie evitate.)

O soluție prezintă într-o serie de mașini moderne este să se dispună de o instrucțiune specială SPECULATIVE\_LOAD care încearcă să citească cuvântul din memoria intermediară, dar dacă nu este acolo, renunță. Dacă valoarea este acolo când este de fapt necesară, se poate folosi, dar dacă nu, hardware-ul trebuie să o aducă imediat. Dacă valoarea nu este de fapt necesară, nu se plătește nici o penalizare pentru ratarea în memoria intermediară.

O situație mult mai rea se poate ilustra cu următoarea instrucțiune:

```
if (x>0) z = y/x;
```

unde x, y și z sunt variabile în virgula mobilă. Să presupunem că toate variabilele sunt citite în registre în avans și împărțirea în virgula mobilă (lentă) este înaltă deasupra testului if. Din nefericire, x este 0 și capcana împărțire-prin-zero termină programul. Rezultatul net este că speculația a cauzat eșecul unui program corect. Mai rău decât atât, programatorul a pus cod explicit pentru a preveni acest lucru și acesta s-a întâmplat.

O posibilă soluție este să se dispună de versiuni speciale de instrucțiuni care ar putea cauza excepții. În plus, un bit, numit bit otrăvitor (poison bit), este adunat la fiecare registru. Când o instrucțiune speculativă specială eșuează, în loc să se genereze o capcană se setează bitul otrăvitor din registrul rezultat. Dacă acest registru este mai târziu accesat de o instrucțiune obișnuită, atunci se activează capcana (așa cum ar trebui). Totuși dacă rezultatul nu este utilizat niciodată, bitul otrăvitor este sters în cele din urmă și nu se întâmplă nimic rău.

## 5.2.8. Microarhitectura UCP Pentium II

Pentium II admite operanți și operații aritmetice pe 32 biți, ca și operații în virgula mobilă pe 64 biți. Admite de asemenea operanți și operații pe 8 și 16 biți, care sunt o mostenire de la procesoarele mai vechi din familie. El poate adresa până la 64 de GB de memorie și citește în memorie cuvinte de 64 biți dintr-o dată. Un sistem tipic Pentium II este ilustrat în fig. 5.34, care este o reprezentare simplificată a procesorului din fig. 5.27.

Așa cum s-a discutat mai devreme, plăcheta Pentium II SEC constă din două circuite integrate: UCP (incluzând memoriile intermediare de nivel 1 separate) și memoria intermediară de nivel 2 unificată. Fig. 5.34 arată componentele primare ale UCP: unitățile Citire/Decodificare, Repartizare/Execuție și Retragere, care lucrează împreună ca o bandă de asamblare de nivel înalt. Aceste trei unități comunică printr-un rezervor de instrucțiuni, un loc pentru păstrarea informațiilor legate de instrucțiuni parțial executate. Informațiile din rezervorul de instrucțiuni sunt memorate într-o tabelă numită ROB (ReOrder Buffer, rom: buffer de reordonare). Pe scurt unitatea de Citire/Decodificare citește instrucțiunile și le sparge în micro-operații pentru memorarea în ROB. Unitatea de Repartizare/Execuție preia micro-operațiile din ROB și le execută. Unitatea Retragere încheie execuția fiecărei micro-operații și actualizează registrele. Instrucțiunile sunt introduse în ROB în ordine, pot fi executate în orice ordine, dar sunt retrase din nou în ordine.

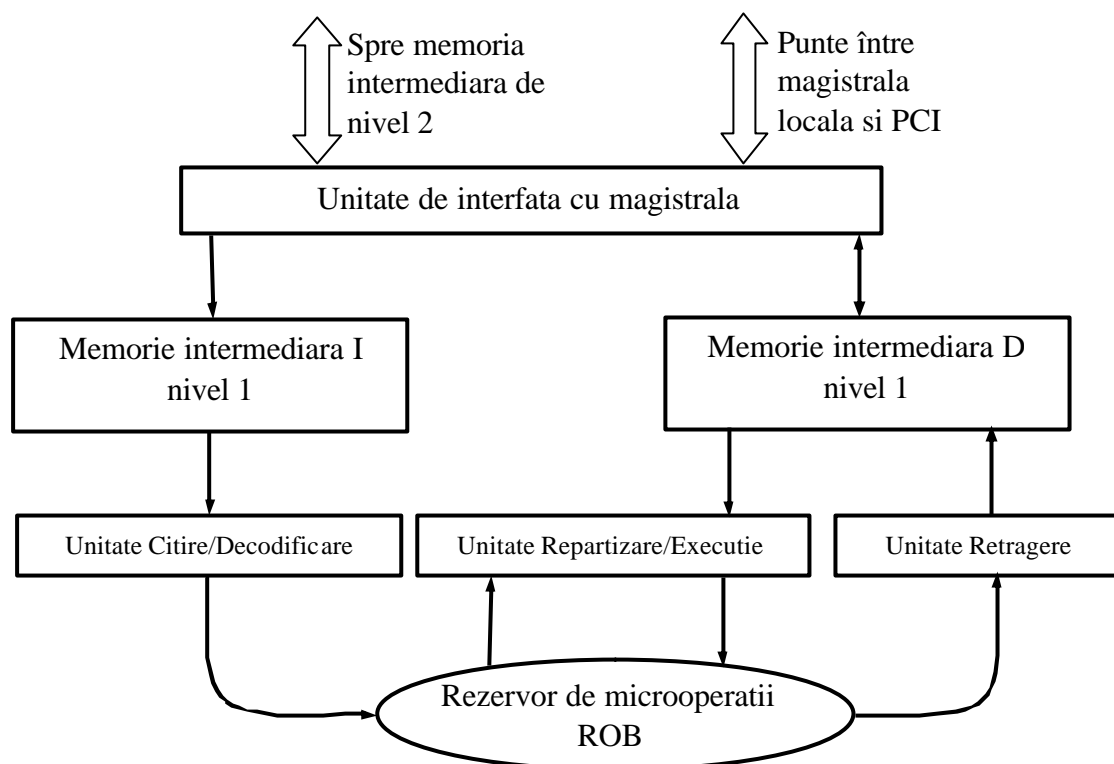


Fig. 5.34. Microarhitectura Pentium II

Unitatea de interfata cu magistrala raspunde de comunicarea cu sistemul de memorie, atât cu memoria intermediara L2 cât si cu memoria principala. Memoria intermediara L2 nu este conectata la magistrala locala, asa ca unitatea de interfata este responsabila cu citirea datelor din memoria principala prin magistrala locala si încarcarea tuturor memoriilor intermediare. Pentium II utilizeaza protocolul MESI de coerența a memoriei intermediare, pe care îl vom descrie când vom ajunge la multiprocesoare.

#### Unitatea Citire/Decodificare

Unitatea Citire/Decodificare este implementata ca banda de asamblare, cu sapte segmente, etichetate de la IFU0 la ROB în Fig. 5.35 (unitatile Repartizare/Executie si Retragere au alte cinci segmente, formând în total o banda de asamblare cu 12 segmente). Instructiunile intra în banda de asamblare în segmentul IFU0, unde sunt încarcate linii întregi de 32 octeti din memoria intermediara I. Ori de câte ori tamponul intern este gol, se copiaza în el o alta linie din memoria intermediara. Registrul NEXT ghideaza procesul de citire.

Deoarece setul de instructiuni Intel, adesea numit IA-32, are instructiuni de lungime variabila, cu mai multe formate, urmatorul segment al benzii de asamblare, IFU1, analizeaza fluxul de octeti pentru a localiza adresa de start a fiecarei instructiuni. Daca este necesar, IFU1 poate cauta în avans pâna la 30 de instructiuni IA-32. Din nefericire, mergând atât de departe înainte, întâlnește patru sau cinci ramificatii conditionate, predictiile nu pot fi corecte pentru fiecare, asa ca valoarea cautarii atât de mult în avans este mica. Segmentul IFU2 aliniaza instructiunile astfel ca urmatorul segment le poate decodifica usor.

Decodificarea începe în segmentul ID0. Decodificarea în Pentium II consta în conversia fiecarei instructiuni IA-2 în una sau mai multe micro-operatii. Instructiunile IA-32 mai simple, ca transferurile registru la registru, se pot converti într-o singura micro-operatie. Altele mai complexe pot necesita pâna la patru micro-operatii. Câteva extrem de complexe necesita chiar mai multe si folosesc ROM-ul secventiatorului de micro-operatii pentru a genera secventa.

Segmentul ID0 are trei decodificatoare interne. Doua din ele sunt pentru instructiuni simple; al treilea le gestioneaza pe celelalte. Ceea ce iese din segmentul ID0 este o secventa de micro-operatii. Fiecare micro-operatie contine un cod de operatie, doua registre sursa si un registru destinatie.

Micro-operatiile sunt plasate în coada în segmentul ID1. Acest segment mai executa si detectia ramificatiilor. Mai întâi se face o predictie statica, pentru orice eventualitate. Predictia depinde de o serie de factori, dar pentru ramificatiile relative la instructiunea curenta se presupune ca vor fi efectuate cele înapoi, iar cele înainte nu vor fi efectuate. Dupa aceasta urmeaza predictorul dinamic de ramificatii, utilizând un algoritm bazat pe istorie (se va vedea când se va trata memoria cache). În loc sa utilizeze numai doi biti pentru istorie el utilizeaza 4 biti. Trebuie sa fie clar ca pentru o banda de asamblare cu 12 segmente penalizarea pentru o eroare

de predicție este enormă, de aceea se utilizează așa de mulți biți de istorie. Dacă ramificația nu este în tabela de istorie, se utilizează predicția statică.

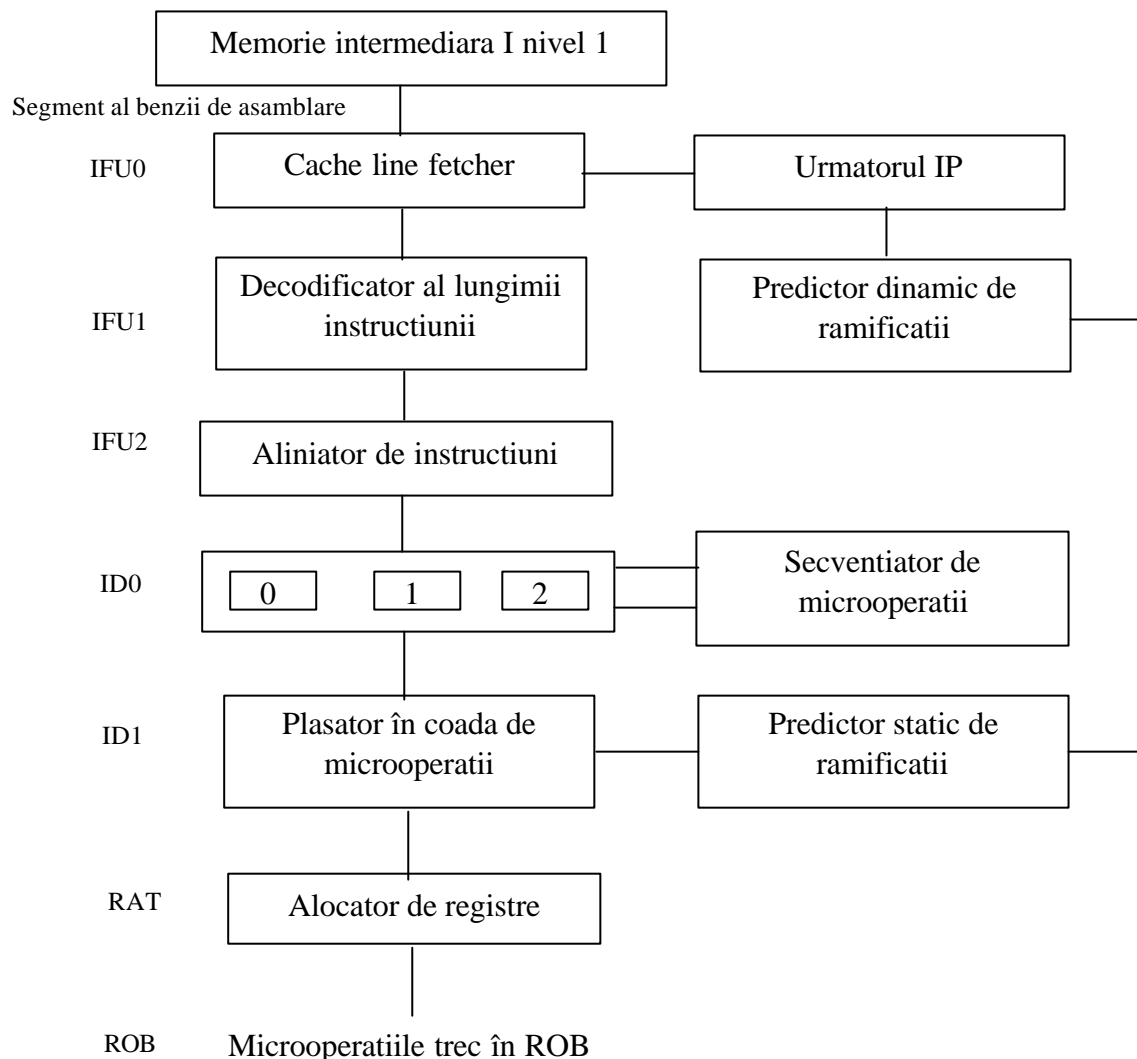


Fig. 5.35. Structura internă a unității Citire/Decodificare (simplificată)

Pentru a evita dependențele WAR și WAW, Pentium II permite redenumirea registrelor, așa cum am văzut în Fig. 5.32. Registrele reale numite în instrucțiunile IA-32 pot fi înlocuite în micro-operatii de oricare din cele 40 de registre temporale interne localizate în ROB. Această redenumire este efectuată în segmentul RAT.

În sfârșit, în ROB sunt depozitate câte trei micro-operatii pe ciclu de ceas. Operanzii sunt de asemenea colectați aici, dacă sunt disponibili. Dacă operanzii unei micro-operatii și registrul rezultat sunt totuși disponibili, iar unitatea de execuție este liberă, aceasta este candidată să fie lansată în execuție. Altfel, rămâne în ROB până când toate resursele sale au fost obținute.

#### Unitatea Repartizare/Execuție

Trecem acum la unitatea Repartizare/Execuție, care este ilustrată în Fig. 5.36. Unitatea de Repartizare/Execuție planifică și execută micro-operatii, rezolvând dependențele și conflictele de resurse. Deși numai trei instrucțiuni ISA pot fi decodificate pe ciclu de ceas (în ID0) pot fi lansate în execuție până la cinci micro-operatii într-un ciclu, câte una pe fiecare port. Aceasta rată nu poate fi susținută pentru că ea depășește capacitatea unității Retragere.

Micro-operatiile pot fi lansate în orice ordine, dar unitatea Retragere le retrage în ordine. Se utilizează o tabelă de scor complexă pentru a gestiona micro-operatii în desfășurare, registre și unități de execuție. Când o micro-operatie este eligibilă pentru execuție ea poate fi începută, chiar dacă altele puse mai devreme în ROB nu



sunt gata. Când mai multe micro-operatii sunt eligibile pentru executie, de catre aceeasi unitate de executie, un algoritm complex o alege pe cea mai buna pentru a fi urmatoarea lansata în executie. De exemplu, executia unei ramificatii este mult mai importanta decât executia unei instructiuni aritmetice, deoarece prima afecteaza fluxul programului.

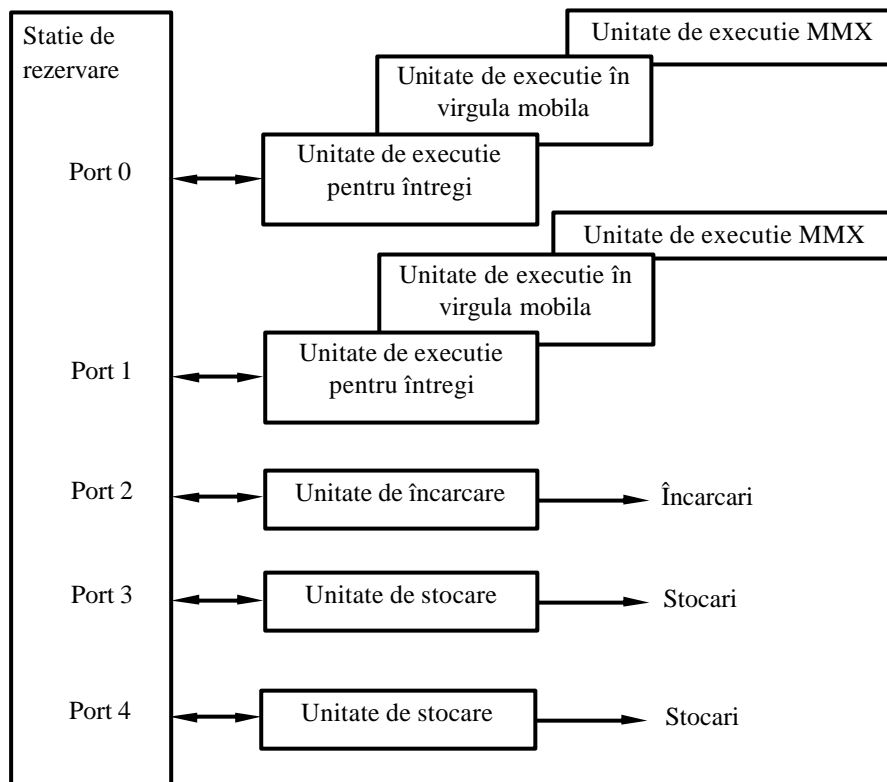


Fig. 5.36. Unitatea de Repartizare/Executie

Unitatea de Repartizare/Executie contine o statie de rezervare si unitati de executie conectate la cinci porturi. Statia de rezervare este o coada cu 20 intrari pentru micro-operatii care au toti operanzii disponibili. Acestea își asteapta rândul în Statia de rezervare pâna când unitatea de executie ceruta devine libera.

Exista cinci porturi de la Statia de rezervare la unitatile de executie. Unele unitati de executie partajeaza un singur port, asa cum s-a aratat. Unitatile de încărcare si Memorie initiaza informatiile corespunzatoare pentru operatiile de încărcare si respectiv memorare. Exista doua porturi pentru memorari. Deoarece numai o singura micro-operatie poate fi lansata per port si per ciclu, daca doua micro-operatii în desfasurare trebuie trimise la acelasi port, una din ele trebuie sa astepte.

### Unitatea de Retragere

O data ce o micro-operatie a fost executata, aceasta merge înapoi la Statia de Rezervare si apoi în ROB pentru a astepta retragerea. Unitatea de Retragere este responsabila pentru trimiterea rezultatelor la locul adecvat - la registrul adecvat, dar si la alte statii din Unitatea Repartizare/Executie asteptând valori. Unitatea de Citire/Decodificare pastreaza registrele "oficiale", adica acele valori de la instructiuni care sau terminat. Unitatea de Retragere pastreaza un set de registre în asteptare, adica acele valori care au fost calculate într-o instructiune care nu s-a terminat pentru ca unele instructiuni precedente nu sunt încă terminate.

Pentium II admite integral executia speculativa, asa ca unele instructiuni vor fi executate inutil si nu vor fi retrase. Aici apare capacitatea de executie înapoi. Daca se constata ca o anumita micro-operatie provine dintr-o instructiune IA-32 care nu ar fi trebuit executata, rezultatele sale sunt eliminate. Este sarcina unitatii de Retragere de a gestiona toate acestea. Pot fi retrase numai instructiuni executate "oficial", iar acestea trebuie retrase în ordinea programului, chiar daca ele au fost poate executate în alta ordine.